



TITLE:

残差と近似解の誤差 (数値計算のアルゴリズムの研究)

AUTHOR(S):

平野, 菅保

CITATION:

平野, 菅保. 残差と近似解の誤差 (数値計算のアルゴリズムの研究). 数理解析研究所講究録 1981, 422: 59-77

ISSUE DATE:

1981-04

URL:

<http://hdl.handle.net/2433/102562>

RIGHT:

残差と近似解の誤差

東芝 平野 菅保

I 有限桁の係数と近似解

実際に電子計算機を使用して数値計算をする場合、 n 元連立1次方程式

$$\sum_{j=1}^n a_{Tij} \cdot x_j - C_{Ti} = 0 \quad i = 1, 2, \dots, n \quad (1)$$

a_{Tij}, C_{Ti} : 実数

の係数 a_{Tij} , 定数 C_{Ti} は、0 を除くと、一般には実数で与えられず、0 でない係数 a_{Tij} , 定数 C_{Ti} は、それぞれ有限桁の数値で与えられるので、それぞれ与えられる最下位の桁の次の桁から下位の桁の数値は、すべて不明である。したがって、不明の桁をすべて 0 として得られる係数 a_{ij} , 定数 C_i を持つ n 元連立1次方程式

$$\sum_{j=1}^n a_{ij} \cdot x_j - C_i = 0 \quad i = 1, 2, \dots, n \quad (2)$$

a_{ij}, C_i : 実数

が n 元連立1次方程式 (1) を代表しているとして、 n 元連立

1 次方程式 (2) を解くことになる。ゆえに、 n 元連立 1 次方程式 (2) の 0 でない係数 a_{ij} , 定数 C_i は

$$a_{Tij} = a_{ij} + \varepsilon_{Tij} \quad -\Delta a_{ij} \leq \varepsilon_{Tij} \leq \Delta a_{ij} \quad (3)$$

$$C_{Ti} = C_i + \varepsilon_{Ti} \quad -\Delta C_i \leq \varepsilon_{Ti} \leq \Delta C_i$$

$$\Delta a_{ij} > 0 \quad \Delta C_i > 0$$

$$\left(\begin{array}{l} \Delta a_{ij} : 0 \text{ でない係数 } a_{ij} \text{ の与えられる最下位の桁の} \\ \text{単位の } 1/2 \\ \Delta C_i : 0 \text{ でない定数 } C_i \text{ の与えられる最下位の桁の単} \\ \text{位の } 1/2 \end{array} \right.$$

を満足することのみ既知である n 元連立 1 次方程式 (1) の係数 a_{Tij} , 定数 C_{Ti} を代表していなければならない。すなわち

$$-\Delta a_{ij} \leq \varepsilon_{ij} \leq \Delta a_{ij} \quad -\Delta C_i \leq \varepsilon_i \leq \Delta C_i \quad (4)$$

を満足する任意の誤差 ε_{ij} , ε_i を実数 a_{ij} , C_i に加えた実数 $(a_{ij} + \varepsilon_{ij})$, $(C_i + \varepsilon_i)$ を係数、定数に持つ n 元連立 1 次方程式

$$\sum_{j=1}^n (a_{ij} + \varepsilon_{ij}) \cdot x_j - (C_i + \varepsilon_i) = 0 \quad (5)$$

$$i = 1, 2, \dots, n$$

の解 $\bar{x}_j (j=1, 2, \dots, n)$ よりも、 n 元連立 1 次方程式 (1) の解に近接していると認めることのできる近似解を、 n 元連立 1 次方程式 (2) から求めることは一般にはできない。したがって、その近似解 $\bar{x}_j (j=1, 2, \dots, n)$ を n 元連立 1 次方程式 (2)

を満足する近似解とする。

II 残差と近似解

n 元連立1次方程式(5)の解 $\bar{x}_j (j=1, 2, \dots, n)$ を n 元連立1次方程式(2)に代入すると、

$$\sum_{j=1}^n a_{ij} \cdot \bar{x}_j - C_i = - \sum_{j=1}^n \varepsilon_{ij} \cdot \bar{x}_j + \varepsilon_i \quad (6)$$

$$i = 1, 2, \dots, n$$

となり、すべての残差(右辺)が0になるとは限らない。すなわち、 n 元連立1次方程式(2)の係数 a_{ij} 、定数 C_i を用いて、数値計算で近似解 $x_j (j=1, 2, \dots, n)$ を求め、その近似解 $x_j (j=1, 2, \dots, n)$ を n 元連立1次方程式(2)に代入して得られる残差が、(4)を満足する適当な誤差 $\varepsilon_{ij}, \varepsilon_i$ を用いて

$$\sum_{j=1}^n a_{ij} \cdot x_j - C_i = - \sum_{j=1}^n \varepsilon_{ij} \cdot x_j + \varepsilon_i \quad (7)$$

$$i = 1, 2, \dots, n$$

とできるならば、近似解 $x_j (j=1, 2, \dots, n)$ は n 元連立1次方程式(2)を満足する近似解である。

解 $x_j (j=1, 2, \dots, n)$ が n 元連立1次方程式(2)を満足しているとの判定をする収束判定に、(4)を用いて、(7)の代りに不等式

$$\left| \sum_{j=1}^n a_{ij} \cdot x_j - C_i \right| \leq \sum_{j=1}^n \Delta a_{ij} \cdot |x_j| + \Delta C_i \quad (8)$$

$$i = 1, 2, \dots, n$$

を用いるが、実際には、必ず(7)を満足する簡便な

$$\left| \sum_{j=1}^n a_{ij} \cdot x_j - C_i \right| \leq \max_{j=1,2,\dots,n} (\Delta a_{ij} \cdot |x_j|, \Delta C_i) \quad (9)$$

$$i = 1, 2, \dots, n$$

を用いる。

Ⅲ 係数行列の数値的特異性

n 元連立1次方程式(2)の係数間に、(4)を満足する適当な ε_{ij} を用いると、

$$\sum_{j=1}^n a_{ij} \cdot \alpha_j = - \sum_{j=1}^n \varepsilon_{ij} \cdot \alpha_j \quad i = 1, 2, \dots, n \quad (10)$$

の関係がなりたつ。少なくとも1つは0でない比例定数 α_j ($j=1, 2, \dots, n$) が存在すれば、 n 元連立1次方程式(2)は数値的に従属な式を含んでいる。

また(10)がなりたつと、

$$\sum_{j=1}^n (a_{ij} + \varepsilon_{ij}) \cdot \alpha_j = 0 \quad i = 1, 2, \dots, n \quad (10')$$

であり、 $(a_{ij} + \varepsilon_{ij})$ による係数行列は正則でないから

$$\sum_{i=1}^n (a_{ij} + \varepsilon_{ij}) \cdot \beta_i = 0 \quad j = 1, 2, \dots, n \quad (11')$$

$$\sum_{i=1}^n a_{ij} \cdot \beta_i = - \sum_{i=1}^n \varepsilon_{ij} \cdot \beta_i \quad j = 1, 2, \dots, n \quad (11)$$

がなりたつ。したがって、 n 元連立1次方程式(2)が数値的に従属な式を含んでいるとの判定をするため、(4)を用いて

、(10)、(11)の代りに

$$\left| \sum_{j=1}^n a_{ij} \cdot \alpha_j \right| \leq \sum_{j=1}^n \Delta a_{ij} \cdot |\alpha_j| \quad i = 1, 2, \dots, n \quad (12)$$

$$\left| \sum_{i=1}^n a_{ij} \cdot \beta_i \right| \leq \sum_{i=1}^n \Delta a_{ij} \cdot |\beta_i| \quad j = 1, 2, \dots, n \quad (13)$$

を用いる。

IV ピボットと係数行列の特異性

ガウスの消去法を用いて n 元連立 1 次方程式 (2) の係数行列 $[a_{ij}]$ を左下三角行列 $[l_{ij}]$ ($i \geq j$) , 右上三角行列 $[m_{ij}]$ ($i \leq j$) に分解すると,

$$\begin{aligned} l_{i1} &= a_{i1} & i &= 1, 2, \dots, n \\ m_{ij} &= a_{ij} / l_{i1} & j &= 1, 2, \dots, n \\ l_{ij} &= a_{ij} - \sum_{k=1}^{\min(i,j)-1} l_{ik} \cdot m_{kj} & \begin{cases} i=2, 3, \dots, n \\ j=2, 3, \dots, i \end{cases} \\ & \quad (i \geq j) \\ m_{ij} &= (a_{ij} - \sum_{k=1}^{\min(i,j)-1} l_{ik} \cdot m_{kj}) / l_{ii} & \begin{cases} i=2, 3, \dots, n \\ j=i, i+1, \dots, n \end{cases} \\ & \quad (i \leq j) \\ (m_{ii} &= 1 \quad i=1, 2, \dots, n) \end{aligned} \quad (14)$$

であるから, n 元連立 1 次方程式 (2) の係数 a_{ij} は

$$a_{ij} = \sum_{k=1}^{\min(i,j)} l_{ik} \cdot m_{kj} \quad (15)$$

である。

(12), (13) を満足する列に関する比例定数 α_j ($j=1, 2, \dots, n$), 行に関する比例定数 β_i ($i=1, 2, \dots, n$) が存在するとして, n 元連立 1 次方程式 (2) の係数行列をスケーリング^{*1} し, 次いで, 完全枢軸選択のガウスの消去法で $(n-1)$ 回消去すると, それぞれ 1 行, 1 列を用いずに α_j ($j=1, 2, \dots, n$), β_i ($i=1, 2, \dots, n$) を求めることになる。(以下では, 用いなかった 1 行, 1 列をそれぞれ第 n 行, 第 n 列とする。)

$$\sum_{j=1}^n a_{ij} \cdot \alpha_j = 0 \quad i = 1, 2, \dots, n \quad (16)$$

$$\sum_{i=1}^n a_{ij} \cdot \beta_i = 0 \quad j = 1, 2, \dots, n \quad (16')$$

($n-1$) 回の枢軸選択によるガウスの前段消去を行なえば

$$\sum_{j=k}^n m_{kj} \cdot \alpha_j = 0 \quad k = 1, 2, \dots, n-1 \quad (17)$$

$$\sum_{i=k}^n l_{ik} \cdot \beta_i = 0 \quad k = 1, 2, \dots, n-1 \quad (17')$$

が得られるので、 α_n, β_n をそれぞれ定数(例えば 1.0)とすれば、(17), (17')から α_j ($j=1, 2, \dots, n-1$), β_i ($i=1, 2, \dots, n-1$) が求められる。したがって、(16), (16')は(17), (17')から

$$\sum_{j=1}^n a_{ij} \cdot \alpha_j = \sum_{k=1}^i l_{ik} \cdot \left(\sum_{j=k}^n m_{kj} \cdot \alpha_j \right) = 0 \quad i = 1, 2, \dots, n-1 \quad (18)$$

$$\sum_{i=1}^n a_{ij} \cdot \beta_i = \sum_{k=1}^j \left(\sum_{i=k}^n l_{ik} \cdot \beta_i \right) \cdot m_{kj} = 0 \quad j = 1, 2, \dots, n-1 \quad (18')$$

$$\sum_{j=1}^n a_{nj} \cdot \alpha_j = \sum_{k=1}^n l_{nk} \cdot \left(\sum_{j=k}^n m_{kj} \cdot \alpha_j \right) = l_{nn} \cdot m_{nn} \cdot \alpha_n \quad (19)$$

$$\sum_{i=1}^n a_{in} \cdot \beta_i = \sum_{k=1}^n \left(\sum_{i=k}^n l_{ik} \cdot \beta_i \right) \cdot m_{kn} = l_{nn} \cdot m_{nn} \cdot \beta_n \quad (19')$$

となる。すなわち、 n 元連立1次方程式(2)の係数行列は、第 n 式を除いて、(12)の条件をすべて満足しているので、

$$|l_{nn} \cdot \alpha_n| \leq \sum_{j=1}^n \Delta a_{nj} \cdot |\alpha_j| \quad m_{nn} = 1.0 \quad (20)$$

を満足していれば、第 n 式($i=n$)も(12)の条件を満足するので、 n 元連立1次方程式(2)は数値的に従属な式を含んでいることになる。

ガウスの消去法で用いるピボットの数値の積の絶対値 $|l_{11} \cdot l_{22} \cdot \dots \cdot l_{nn}|$ は、ピボットの数値の採用の仕方に関係なく、係数行列の行列式の値の絶対値に等しく、一方、完全枢軸選択のガウスの消去法では、各消去段階でピボット l_{ii} ($i = 1, 2, \dots, n$) として採用してよい係数の中で、絶対値最大の数値をピボットとして採用しているので、最後 (n 回目) のピボット l_{nn} の絶対値は他のピボット $l_{11}, l_{22}, \dots, l_{n-1, n-1}$ の絶対値に比較して、一般に小さい。特に係数行列が悪条件の時にそれが著しい。

また n 元連立 1 次方程式 (2) の係数行列をスケーリング^{*1} し、次いで、完全枢軸選択のガウスの消去法で消去すると、 $\alpha_n = 1.0$, $\beta_n = 1.0$ ならば

$$|\alpha_j| \leq 1.0 \quad |\beta_i| \leq 1.0 \quad i, j = 1, 2, \dots, n-1$$

がほぼなりたつ。

*1 n 元連立 1 次方程式 (2) の係数の絶対値が極端に異なる時は、次のようなスケーリングを行なう。

$$A_i = \max_{j=1, 2, \dots, n} (|a_{ij}|) \quad i = 1, 2, \dots, n$$

$$\bar{A}_j = \max_{i=1, 2, \dots, n} (|a_{ij}| / A_i) \quad j = 1, 2, \dots, n$$

$$a_{nij} = a_{ij} / (A_i \cdot \bar{A}_j) \quad i, j = 1, 2, \dots, n$$

新しい係数行列 a_{nij} の列に関する比例定数 α_j ($j = 1, 2, \dots, n$)、行に関する比例定数 β_i ($i = 1, 2, \dots, n$) から

$$\alpha_j = \alpha_{Nj} / \bar{A}_j \quad j = 1, 2, \dots, n$$

$$\beta_i = \beta_{Ni} / A_i \quad i = 1, 2, \dots, n$$

によって n 元連立 1 次方程式 (2) の係数行列の列に関する比例定数 α_j ($j = 1, 2, \dots, n$), 行に関する比例定数 β_i ($i = 1, 2, \dots, n$) を求める。

条件 (4) を満足する, すべての ε_{ij} ($i, j = 1, 2, \dots, n$) を考慮しているわけではないので, 条件 (20) を満足しなくても, 条件 (12) を満足する場合がある。

そこで, あまり実用にはならないが, (12) を満足しない場合を求めてみる。まず $\alpha_n = 1.0$ と (17) から α_j ($j = 1, 2, \dots, n$) を求め, 次のようなスケーリングをすると, n 元連立 1 次方程式 (2) の係数は

$$\sum_{j=1}^n \bar{a}_{ij} \cdot \bar{\alpha}_j = 0 \quad i = 1, 2, \dots, n \quad (21)$$

$$\bar{a}_{ij} = (a_{ij} \cdot \alpha_j) / \left(\sum_{k=1}^n \Delta a_{ik} \cdot |\alpha_k| \right)$$

$$\Delta \bar{a}_{ij} = (\Delta a_{ij} \cdot \alpha_j) / \left(\sum_{k=1}^n \Delta a_{ik} \cdot |\alpha_k| \right)$$

となり, n 元連立 1 次方程式 (2) が数値的に従属な式を含んでいるとの判定をするため, (12) の代りに

$$\left| \sum_{j=1}^n \bar{a}_{ij} \cdot \bar{\alpha}_j \right| \leq \sum_{j=1}^n \Delta \bar{a}_{ij} \cdot |\bar{\alpha}_j| \quad (22)$$

$$i = 1, 2, \dots, n$$

を用いる。 $\bar{\alpha}_n = 1.0$ とすると,

$$\bar{\alpha}_j \doteq 1.0 \quad j = 1, 2, \dots, n-1$$

であり、ほぼ次式がなりたつ。

$$\left| \sum_{j=1}^n \bar{a}_{ij} \cdot \bar{\alpha}_j \right| \leq 1.0 \quad i=1, 2, \dots, n \quad (24)$$

次いで、完全枢軸選択のガウスの消去法で係数行列の \bar{a}_{ij} を \bar{l}_{ij} ($i \geq j$), \bar{m}_{ij} ($i \leq j$) に分解し、(17), (18), (19) と同様に

$$\sum_{j=k}^n \bar{m}_{kj} \cdot \bar{\alpha}_j = 0 \quad k=1, 2, \dots, n-1 \quad (25)$$

$$\sum_{j=1}^n \bar{a}_{ij} \cdot \bar{\alpha}_j = \sum_{k=1}^i \bar{l}_{ik} \cdot \left(\sum_{j=k}^n \bar{m}_{kj} \cdot \bar{\alpha}_j \right) = 0 \quad (26)$$

$$i=1, 2, \dots, n-1$$

$$\sum_{j=1}^n \bar{a}_{nj} \cdot \bar{\alpha}_j = \sum_{k=1}^n \bar{l}_{nk} \cdot \left(\sum_{j=k}^n \bar{m}_{kj} \cdot \bar{\alpha}_j \right) = \bar{l}_{nn} \cdot \bar{m}_{nn} \cdot \bar{\alpha}_n \quad (27)$$

$$(\bar{m}_{nn} = 1.0)$$

となる。n元連立1次方程式(21)の第n式を除いて、(22)の条件をすべて満足しているので

$$|\bar{l}_{nn} \cdot \bar{\alpha}_n| \leq \sum_{j=1}^n \Delta \bar{a}_{nj} \cdot |\bar{\alpha}_j| \quad (28)$$

を満足していれば、n元連立1次方程式(2)の第n式も(20)の条件を満足するので、n元連立1次方程式(2)は数値的に従属な式を含んでいることになる。

n元連立1次方程式(21)の第n式を除いて、(26)では、残差を0としているが、第i式($i \neq n$)の残差を(22)の右辺

$$\sum_{j=1}^n \Delta \bar{a}_{ij} \cdot |\bar{\alpha}_j| \quad (\alpha_n = 1.0) \quad (29)$$

$$i=1, 2, \dots, n-1$$

とすると、比例定数 $\bar{\alpha}_j$ ($j=1, 2, \dots, n-1$) の変化分は

$$b_{ji} \cdot \left(\sum_{k=1}^n \Delta \bar{a}_{ik} \cdot |\bar{\alpha}_k| \right) \quad (30)$$

$i, j = 1, 2, \dots, n-1$

(ただし、係数行列 $[b_{ij}]$ は係数行列 $[a_{ij}]$ ($i, j = 1, 2, \dots, n-1$) の逆行列)

となる。したがって、(22)の第1式から第($n-1$)式までの右辺を、それぞれ残差とすると、比例定数 $\bar{\alpha}_j$ ($j = 1, 2, \dots, n-1$) の変化分 $\Delta \bar{\alpha}_j$ ($j = 1, 2, \dots, n-1$) は

$$\Delta \bar{\alpha}_j = \sum_{i=1}^{n-1} b_{ji} \cdot \left(\sum_{k=1}^n \Delta \bar{a}_{ik} \cdot |\bar{\alpha}_k| \right) \quad (31)$$

$j = 1, 2, \dots, n-1$

$$\Delta \bar{\alpha}_n = 0, \quad \bar{\alpha}_n = 1.0$$

となる。ここで $\bar{\alpha}_j + \Delta \bar{\alpha}_j$ ($j = 1, 2, \dots, n$) を n 元連立1次方程式(21)の第 n 式に代入する。

$$\begin{aligned} & \sum_{j=1}^n \bar{a}_{nj} \cdot (\bar{\alpha}_j + \Delta \bar{\alpha}_j) \\ &= \bar{l}_{nn} \cdot \bar{m}_{nn} \cdot \bar{\alpha}_n + \sum_{j=1}^{n-1} \bar{a}_{nj} \cdot \left\{ \sum_{i=1}^{n-1} b_{ji} \cdot \left(\sum_{k=1}^n \Delta \bar{a}_{ik} \cdot |\bar{\alpha}_k| \right) \right\} \\ &= \bar{l}_{nn} + \sum_{i=1}^{n-1} \left\{ \left(\sum_{j=1}^{n-1} \bar{a}_{nj} \cdot b_{ji} \right) \cdot \left(\sum_{k=1}^n \Delta \bar{a}_{ik} \cdot |\bar{\alpha}_k| \right) \right\} \quad (32) \\ & \bar{m}_{nn} = 1.0, \quad \bar{\alpha}_n = 1.0 \end{aligned}$$

(32)より、次式を満足すれば、 n 元連立1次方程式(2)は数値的に従属な式を含まないとみなしている。

$$\begin{aligned} |\bar{l}_{nn} \cdot \bar{m}_{nn} \cdot \bar{\alpha}_n| &> \sum_{i=1}^{n-1} \left\{ \left(\left| \sum_{j=1}^{n-1} \bar{a}_{nj} \cdot b_{ji} \right| \right) \cdot \left(\sum_{k=1}^n \Delta \bar{a}_{ik} \cdot |\bar{\alpha}_k| \right) \right\} \\ &+ \sum_{j=1}^n \Delta \bar{a}_{nj} \cdot |\bar{\alpha}_j + \Delta \bar{\alpha}_j| \quad (33) \end{aligned}$$

ただし、(33)の中に含まれる $\Delta \bar{\alpha}_j$ は次式で示される。

$$\Delta \bar{\alpha}_j = \sum_{i=1}^{n-1} \bar{b}_{ji} \cdot \left(\sum_{k=1}^n \Delta \bar{a}_{ik} \cdot |\bar{\alpha}_k| \right) \quad j=1, 2, \dots, n-1 \quad (34)$$

$$\bar{b}_{ji} = b_{ji}, \quad \left(\sum_{k=1}^{n-1} \bar{a}_{nk} \cdot b_{ki} \right) \geq 0$$

$$\bar{b}_{ji} = -b_{ji}, \quad \left(\sum_{k=1}^{n-1} \bar{a}_{nk} \cdot b_{ki} \right) < 0$$

しかし、比例定数 $\bar{\alpha}_j$ が $(\bar{\alpha}_j + \Delta \bar{\alpha}_j)$ に変化したので、条件 (22) は

$$\left| \sum_{j=1}^n \bar{a}_{ij} \cdot (\bar{\alpha}_j + \Delta \bar{\alpha}_j) \right| \leq \sum_{j=1}^n \Delta \bar{a}_{ij} \cdot |\bar{\alpha}_j + \Delta \bar{\alpha}_j| \quad (35)$$

$i=1, 2, \dots, n$

となる。したがって、(33) の条件で n 元連立 1 次方程式 (2) が数値的に従属な式を含まないとは、厳密にはいえないが、本来の誤差の性質から考えて、実用的には充分である。

ここで n 元連立 1 次方程式 (2) の誤差を含む係数行列の条件数として

$$\rho = \min \frac{\left| \sum_{i=1}^n \beta_i \cdot \left(\sum_{j=1}^n a_{ij} \cdot \alpha_j \right) \right|}{\sum_{i=1}^n |\beta_i| \cdot \left(\sum_{j=1}^n \Delta a_{ij} \cdot |\alpha_j| \right)} \quad (36)$$

ただし、係数行列 $[a_{ij}]$ は正則であり、且つ比例定数 α_j, β_i ($i, j=1, 2, \dots, n$) はそれぞれ少なくとも 1 つは 0 でない。

を定義する。この条件数が

$$\rho \gg 1.0$$

であれば、誤差を含む係数行列の性質が良く、 ρ が 1 に近ければ近い程、この係数行列を持つ n 元連立 1 次方程式 (2) は

式の間、の独立性を数値的にだんだん失ってゆく。

V 枢軸 (ピボット) と固有値

n 元連立 1 次方程式 (2) の係数をスケールリング^{*1} し、次いで、完全枢軸選択によるガウスの消去法で計算したとき、計算途中において枢軸 (ピボット) の数値に桁落ち誤差が入る場合、 n 元連立 1 次方程式 (2) の近似解を求めるとともに、 n 元連立 1 次方程式 (2) の係数行列 $[a_{ij}]$ をスケールリング^{*1} して得られる係数行列の列および行の比例定数 $\alpha_j^{(0)}, \beta_i^{(0)}$ ($i, j = 1, 2, \dots, n$) を (16), (16) を用いて求める。

この比例定数 $\alpha_j^{(0)}, \beta_i^{(0)}$ ($i, j = 1, 2, \dots, n$) とスケールリング^{*1} して得られる係数行列とで

$$\left| \sum_{j=1}^n a_{ij}^{(0)} \cdot \alpha_j^{(0)} \right| \quad i = 1, 2, \dots, n \quad (37)$$

$$\left| \sum_{i=1}^n a_{ij}^{(0)} \cdot \beta_i^{(0)} \right| \quad j = 1, 2, \dots, n \quad (37')$$

(ただし、 $a_{ij}^{(0)}$ は n 元連立 1 次方程式 (2) の係数行列をスケールリング^{*1} して得られる係数行列

を計算すれば、(18), (18'), (19), (19') と枢軸 (ピボット) が桁落ちしていることから

$$\left| \sum_{j=1}^n a_{ij}^{(0)} \cdot \alpha_j^{(0)} \right| \doteq 0 \quad i = 1, 2, \dots, n-1 \quad (38)$$

$$\left| \sum_{i=1}^n a_{ij}^{(0)} \cdot \beta_i^{(0)} \right| \doteq 0 \quad j = 1, 2, \dots, n-1 \quad (38')$$

$$\left| \sum_{j=1}^n a_{nj}^{(0)} \cdot \alpha_j^{(0)} \right| = \left| l_{nn}^{(0)} \cdot m_{nn}^{(0)} \cdot \alpha_n^{(0)} \right| \ll \max_{j=1, 2, \dots, n} \left(\left| a_{nj}^{(0)} \cdot \alpha_j^{(0)} \right| \right) \quad (39)$$

$$\left| \sum_{i=1}^n a_{in}^{(0)} \cdot \beta_i^{(0)} \right| = \left| l_{nn}^{(0)} \cdot m_{nn}^{(0)} \cdot \beta_n^{(0)} \right| \ll \max_{i=1, 2, \dots, n} \left(\left| a_{in}^{(0)} \cdot \beta_i^{(0)} \right| \right) \quad (39')$$

となる。すなわち、この比例定数 $\alpha_j^{(0)}, \beta_i^{(0)}$ ($i, j = 1, 2, \dots, n$) は、一般に、係数行列 $[a_{ij}^{(0)}]$ の含む固有値の中で絶対値最小の固有値の固有ベクトルの近似ベクトルである。したがってこの場合、 n 元連立 1 次方程式 (2) の係数行列 $[a_{ij}]$ をスケールング^{*1}して得られる係数行列 $[a_{ij}^{(0)}]$ の絶対値最大の固有値 $\lambda_{\max}^{(0)}$ と絶対値最小の固有値 $\lambda_{\min}^{(0)}$ との関係が

$$|\lambda_{\max}^{(0)} / \lambda_{\min}^{(0)}| \gg 1.0 \quad (40)$$

である。しかし、スケールング^{*1}されていない一般の係数行列の場合は、(40) を満足しているからといって、(38), (38'), (39), (39') を満足するとは限らない。

比例定数 $\alpha_j^{(0)}, \beta_i^{(0)}$ ($i, j = 1, 2, \dots, n$) を対応する固有ベクトルの近似ベクトルとした場合の精度の桁数は、(38), (38'), (39), (39') から、それぞれ

$$\log_e \left\{ \max_{j=1,2,\dots,n} (|a_{nj}^{(0)} \cdot \alpha_j^{(0)}|) / |l_{nn}^{(0)} \cdot m_{nn}^{(0)} \cdot \alpha_n^{(0)}| \right\} \quad (41)$$

$$\log_e \left\{ \max_{i=1,2,\dots,n} (|a_{in}^{(0)} \cdot \beta_i^{(0)}|) / |l_{nn}^{(0)} \cdot m_{nn}^{(0)} \cdot \beta_n^{(0)}| \right\} \quad (41')$$

ただし、 m 進法

の程度である。

次に、 n 元連立 1 次方程式 (2) の係数行列 $[a_{ij}]$ をスケールング^{*1}して得られる係数行列 $[a_{ij}^{(0)}]$ の比例定数 $\alpha_i^{(0)}, \beta_i^{(0)}$ ($i, j = 1, 2, \dots, n$) に対応する固有値、固有ベクトルをそれぞれ $\lambda_T^{(0)}, \alpha_{Tj}^{(0)}, \beta_{Ti}^{(0)}$ ($i, j = 1, 2, \dots, n$) とすると

$$\sum_{i=1}^n \beta_{Ti}^{(0)} \cdot \left(\sum_{j=1}^n a_{ij}^{(0)} \cdot \alpha_{Tj}^{(0)} \right) = \lambda_T^{(0)} \cdot \sum_{i=1}^n \beta_{Ti}^{(0)} \cdot \alpha_{Ti}^{(0)} \quad (42)$$

であり、この計算での桁落ち誤差の入る桁数はほぼ

$$P_T = \log_t \left\{ \frac{\sum_{i=1}^n |\beta_{Ti}^{(o)}| \cdot \left(\sum_{j=1}^n |a_{ij}^{(o)} \cdot \alpha_{Tj}^{(o)}| \right)}{\left| \lambda_T^{(o)} \cdot \sum_{i=1}^n \beta_{Ti}^{(o)} \cdot \alpha_{Ti}^{(o)} \right|} \right\} \quad (43)$$

$$P_T > 1 \quad t: t \text{進数}$$

である。また、 $\alpha_j^{(o)}, \beta_i^{(o)} (i, j=1, 2, \dots, n)$ の誤差の入っていない術数をそれぞれ

$$P_\alpha = \min_{j=1, 2, \dots, n} \log_t \left\{ \left| \alpha_{Tj}^{(o)} / \Delta \alpha_{Tj}^{(o)} \right| \right\}, \quad \alpha_{Tj}^{(o)} = \alpha_j^{(o)} + \Delta \alpha_{Tj}^{(o)} \\ P_\alpha > 1 \quad (44)$$

$$P_\beta = \min_{i=1, 2, \dots, n} \log_t \left\{ \left| \beta_{Ti}^{(o)} / \Delta \beta_{Ti}^{(o)} \right| \right\}, \quad \beta_{Ti}^{(o)} = \beta_i^{(o)} + \Delta \beta_{Ti}^{(o)} \\ P_\beta > 1 \quad (45)$$

とすると

$$\sum_{i=1}^n \beta_i^{(o)} \cdot \left(\sum_{j=1}^n a_{ij}^{(o)} \cdot \alpha_j^{(o)} \right) = \sum_{i=1}^n \beta_{Ti}^{(o)} \cdot \left(\sum_{j=1}^n a_{ij}^{(o)} \cdot \alpha_{Tj}^{(o)} \right) - \lambda_T^{(o)} \cdot \sum_{j=1}^n \beta_{Tj}^{(o)} \cdot \Delta \alpha_{Tj}^{(o)} \\ - \lambda_T^{(o)} \cdot \sum_{i=1}^n \Delta \beta_{Ti}^{(o)} \cdot \alpha_{Ti}^{(o)} + \sum_{i=1}^n \Delta \beta_{Ti}^{(o)} \cdot \left(\sum_{j=1}^n a_{ij}^{(o)} \cdot \Delta \alpha_{Tj}^{(o)} \right) \\ \sum_{i=1}^n \beta_{Ti}^{(o)} \cdot \left(\sum_{j=1}^n a_{ij}^{(o)} \cdot \alpha_{Tj}^{(o)} \right) = \lambda_T^{(o)} \cdot \sum_{i=1}^n \beta_{Ti}^{(o)} \cdot \alpha_{Ti}^{(o)} \quad (46)$$

の式の中では、(44), (45) から

$$\left| \lambda_T^{(o)} \cdot \sum_{i=1}^n \beta_{Ti}^{(o)} \cdot \alpha_{Ti}^{(o)} \right| \gg \left| \lambda_T^{(o)} \cdot \sum_{j=1}^n \beta_{Tj}^{(o)} \cdot \Delta \alpha_{Tj}^{(o)} \right| \\ \left| \lambda_T^{(o)} \cdot \sum_{i=1}^n \beta_{Ti}^{(o)} \cdot \alpha_{Ti}^{(o)} \right| \gg \left| \lambda_T^{(o)} \cdot \sum_{i=1}^n \Delta \beta_{Ti}^{(o)} \cdot \alpha_{Ti}^{(o)} \right| \quad (47)$$

なる関係があり

$$\left| \lambda_T^{(o)} \cdot \sum_{i=1}^n \beta_{Ti}^{(o)} \cdot \alpha_{Ti}^{(o)} \right| \gg \left| \sum_{i=1}^n \Delta \beta_{Ti}^{(o)} \cdot \left(\sum_{j=1}^n a_{ij}^{(o)} \cdot \Delta \alpha_{Tj}^{(o)} \right) \right| \quad (48)$$

を満足するならば、すなわち、(43), (44), (45) から

$$P_T < P_\alpha + P_\beta \quad (49)$$

であるならば、

$$\sum_{i=1}^n \beta_i^{(0)} \cdot \left(\sum_{j=1}^n a_{ij}^{(0)} \cdot \alpha_j^{(0)} \right) \doteq \lambda_T^{(0)} \cdot \sum_{i=1}^n \beta_{Ti}^{(0)} \cdot \alpha_{Ti}^{(0)} \quad (50)$$

である。したがって、固有値は

$$\begin{aligned} \lambda_T^{(0)} &\doteq \sum_{i=1}^n \beta_i^{(0)} \cdot \left(\sum_{j=1}^n a_{ij}^{(0)} \cdot \alpha_j^{(0)} \right) / \left(\sum_{i=1}^n \beta_{Ti}^{(0)} \cdot \alpha_{Ti}^{(0)} \right) \\ &\doteq l_{nn}^{(0)} \cdot m_{nn}^{(0)} \cdot \beta_n^{(0)} \cdot \alpha_n^{(0)} / \left(\sum_{i=1}^n \beta_i^{(0)} \cdot \alpha_i^{(0)} \right) \end{aligned} \quad (51)$$

で求められる。

ここで、係数行列 $[a_{ij}^{(0)}]$ の最後の列を

$$a_{in}^{(0)} \longrightarrow 10^{-m} \cdot a_{in}^{(0)} \quad m \gg 1 \quad (52)$$

$i = 1, 2, \dots, n$

とすると、固有値は

$$\lambda_T^{(0)} \doteq \frac{\sum_{i=1}^n \beta_i^{(0)} \cdot \left\{ \sum_{j=1}^{n-1} a_{ij}^{(0)} \cdot \alpha_j^{(0)} + (10^{-m} \cdot a_{in}^{(0)}) \cdot (10^m \cdot \alpha_n^{(0)}) \right\}}{\sum_{i=1}^{n-1} \beta_i^{(0)} \cdot \alpha_i^{(0)} + \beta_n^{(0)} \cdot (10^m \cdot \alpha_n^{(0)})} \quad (53)$$

となり、 λ_T の絶対値は小さくなる。しかし、係数行列の条件数(36)は変化したわけではない。すなわち、(40)の条件が満足されても、係数行列の性質が必ずしも悪いわけではなく、(40)の条件を満足している場合、係数行列の性質が悪いとほいえるのは、スケーリング^{*1}された係数行列の場合である

また、次式を満足する固有値 λ_T は誤差のみである。

$$|\lambda_T \cdot \alpha_{Ti}| \leq \left| \sum_{j=1}^n \Delta a_{ij} \cdot \alpha_{Tj} \right|, \quad i = 1, 2, \dots, n \quad (54)$$

ただし、 λ_T は n 元連立 1 次方程式 (2) の係数行列の固有値であり、 α_{Ti} ($i = 1, 2, \dots, n$) は固有値 λ_T に対応

する固有ベクトルである。

VI 係数および定数の誤差と近似解の誤差

n 元連立1次方程式(2)の係数行列 $[a_{ij}]$ の逆行列を

$$A^{-1} = \bar{B} = \begin{bmatrix} \bar{b}_{ij} \end{bmatrix} \quad (55)$$

とすると、 n 元連立1次方程式(2)の解 x_j ($j=1, 2, \dots, n$)の誤差は、ほぼ

$$\begin{aligned} \Delta x_{Tj} &\doteq \sum_{i=1}^n \bar{b}_{ji} \cdot \delta_{Ti} & j=1, 2, \dots, n \\ \delta_{Ti} &= -\sum_{j=1}^n \varepsilon_{Tij} \cdot x_j + \varepsilon_{Ti} & i=1, 2, \dots, n \end{aligned} \quad (56)$$

である。したがって

$$\begin{aligned} -\sum_{i=1}^n |\bar{b}_{ji}| \cdot \Delta d_{Ti} &\leq \Delta x_{Tj} \leq \sum_{i=1}^n |\bar{b}_{ji}| \cdot \Delta d_{Ti} \\ \Delta d_{Ti} &= \sum_{j=1}^n \Delta a_{ij} \cdot |x_j| + \Delta C_i & j=1, 2, \dots, n \end{aligned} \quad (57)$$

であるが、実用としては

$$\begin{aligned} -\max_{i=1, 2, \dots, n} (|\bar{b}_{ji}| \cdot \Delta d_i) &\leq \Delta x_j \leq \max_{i=1, 2, \dots, n} (|\bar{b}_{ji}| \cdot \Delta d_i) \\ \Delta d_i &= \max_{j=1, 2, \dots, n} (\Delta a_{ij} \cdot |x_j|, \Delta C_i) & j=1, 2, \dots, n \end{aligned} \quad (58)$$

を用いている。

また、(56)からわかるように、解 x_j ($j=1, 2, \dots, n$)の誤差 Δx_{Tj} の範囲を(57), (58)で示したが、すべての解が同時に最大の誤差を含むわけではない。

VII 補足

今回は、 n 元連立1次方程式(2)の係数 a_{ij} 、定数 C_i が含む誤差 ε_{ij} 、 ε_i はそれぞれ独立であるとして考えてきたが、一般に、独立でないときには、今回の解の誤差より解の誤差が小さくなる。また、誤差 ε_{ij} 、 ε_i が独立のときの例は参考文献①、②、③、④にある。

n 元連立1次方程式(2)を完全枢軸選択によるガウスの消去法で計算する場合、必要以上の計算桁数をとらないで、

$$\log_t \left\{ \max_{i,j=1,2,\dots,n} (|a_{ij} / \Delta a_{ij}|) \right\} \quad (59)$$

t : t 進数

に、丸めの誤差を防ぐために必要な桁を数桁加えた桁数で計算するためには、スケーリング^{*1}の方法として

$$a_{ij}^{(0)} = (a_{ij} \cdot x_j) / \max_{j=1,2,\dots,n} (|\Delta a_{ij} \cdot x_j|, \Delta C_i) \quad (60)$$

$$i, j = 1, 2, \dots, n$$

(ただし、 x_j ($j=1, 2, \dots, n$) は n 元連立1次方程式(2)の解

が好ましい。しかし、 x_j ($j=1, 2, \dots, n$) の値は未知であるので、近似解が既知の場合には、その近似解^{*2}を用いる。

*2 (60)の分母で選ばれる項

$$\Delta a_{im} \cdot x_{im} = \max_{j=1,2,\dots,n} (|\Delta a_{ij} \cdot x_j|, \Delta C_i) \quad (61)$$

$$i = 1, 2, \dots, n$$

に含まれる解 x_{im} ($i=1, 2, \dots, n$) およびピボットと

して採用される係数を持つ解は、少なくとも1～2桁の精度を必要とする。

最後に、ピボットに関することであるが、スケーリング^{*1}をした後、完全枢軸選択によるガウスの消去法を用いる場合、ピボットの積 $|l_{11}^{(0)} \cdot l_{22}^{(0)} \cdots l_{nn}^{(0)}|$ が係数行列の行列式の値の絶対値に等しく、ピボットの選び方とは関係なく一定であるから、ピボットとして用いてよい係数の中で絶対値最大の数値を用いてゆくと、前段消去の最後のほうで採用するピボットの数値の絶対値が小さくなり、誤差が入る。したがって、「ピボットとして用いてよい係数の中で絶対値最大の数値をピボットに用いず、絶対値の大きい順に数えて数番目の数値をピボットとして採用することが好ましい。」との考えもあるが、この考えは、 n 元連立1次方程式(2)の係数 a_{ij} に誤差を考え、数値的に従属な式を含むか否かを考えれば、今回述べたことから、意味がないことがわらう。

参考文献

①「数値の誤差と数式の誤差」 平野菅保

日本物理学会 応用数学力学講演会(1965年)

②「代数方程式の解法及び誤差」 平野菅保

情報処理学会

第8回プログラミング・シンポジウム(1967年)

- ③ 「関数近似と誤差の問題」 平野管保
電子計算機のための数値計算法Ⅱ 培風館
山内二郎、森口繁一、一松信 共編 (1967年)
- ④ 「INPUT の精度(けた数)で行う計算」 平野管保
技術者のための電子演算手法 コロナ社
木村久男編 (1970年)
- ⑤ 「多元連立1次方程式の精密な解き方」 平野管保
コンピューターによる構造工学講座Ⅱ-1-A
戸川隼人、藤井宏、三好哲彦、平野管保 共著
培風館 (1971年)
- ⑥ 「連立方程式の数値解の一つの試み」 平野管保
数理解析研究所講究録107
京都大学数理解析研究所 (1970年)